

## INTRODUCTION

Due to the uncultivated status of the majority of microorganisms in nature, little is known about their genetic properties, biochemical functions, and metabolic characteristics. Although sequence determination of the microbial community 'genome' is now possible with high throughput sequencing technology, the complexity and magnitude of most microbial communities make meaningful data acquisition and interpretation difficult. Therefore, we are sequencing groundwater microbial communities with manageable diversity and complexity (~10-400 phylogenets) at the U.S. Department of Energy's Environmental Remediation Sciences Program (ERSP) Field Research Center (FRC), Oak Ridge, TN. The microbial community has been sequenced from a groundwater sample contaminated with very high levels of nitrate, uranium and other heavy metals and pH ~3.7. Sequence analysis of this groundwater sample based on a 16S rDNA library revealed 10 operational taxonomic units (OTUs) at the 99.6% cutoff with >90% of the OTUs represented by an unidentified  $\gamma$ -proteobacterial species similar to *Frateruia*. Additional OTUs were related to a  $\beta$ -proteobacterial species of the genus *Azarcus*. Three clone libraries with different DNA fragment sizes (3, 8 and 40 kb) were constructed, and 50-60 Mb raw sequences were obtained using a shotgun sequencing approach. The raw sequences were assembled into 2770 contigs totaling ~6 Mb which were further assembled into 224 scaffolds (1.8 kb-2.4 Mb). Preliminary binning of the scaffolds suggest 4 primary groupings (2 *Frateruia*-like  $\gamma$ -proteobacteria, 1 *Burkholderia*-like  $\beta$ -proteobacteria and 1 *Herbaspirillum*-like  $\beta$ -proteobacteria). Genes identified from the sequences were consistent with the geochemistry of the site, including multiple nitrate reductase and metal resistance genes. A low level of strain diversity was observed in the sample, with little significant polymorphism detected in the ORFs studied. We hypothesize that the major adaptive response within the community following site contamination resulted from lateral gene transfer events within the community followed by adaptive evolution of individual genetic elements. These adaptive events likely triggered multiple selective sweeps within the populations that have reduced the strain heterogeneity of the community.



Fig. 1. Map showing the location of Area 3, well FW106 of the ERSF FRC site from which the microbial community genome sequences were obtained.

## FW106 Groundwater Geochemistry

- Uranium – 51 mg/L (soil ~500 mg/kg)
- Nitrate – 2,331 mg/L
- Sulfate – 1997 mg/L
- Total Organic Carbon (TOC) – 244 mg/L
- Total Inorganic Carbon (TIC) – 284 mg/L
- pH – 3.7
- Heavy Metals and Organics –  $\mu$ g to mg/L

## Metagenome Statistics

- ~70 Mb raw sequence
- ~8 Mb assembled sequence
- 2770 contigs
- 224 scaffolds (1.8 kb-2.4 Mb)
- 5 preliminary phylogenies based on identification of anchor genes (16S rRNA, 23S rRNA, *gyrB*, *recA*, *rpoB*, *ileS*, *fusA*)
- ~70 of metagenomic sequences assigned to these 5 phylogenies
- Estimated Genome Coverage<sup>a</sup>
  - FRC Gamma Group I (*Frateruia* I) – 9.7X
  - FRC Beta Group I (*Burkholderia*) – 1.1X
  - FRC Gamma Group II (*Frateruia* II) – <0.1X
  - FRC Beta Group II (*Herbaspirillum*) – <0.1X
  - FRC Alpha Group I (*Afrapia*) – <0.1X

<sup>a</sup> Genome coverage estimated using reference genome sizes as follows: *Burkholderia*, 7.8 Mb; *Herbaspirillum*, 4.8 Mb; *Frateruia*, 4.2 Mb; *Herbaspirillum*, 4.1 Mb.

## Significant Scaffolds

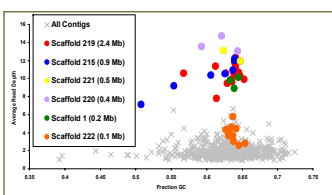


Fig. 2. Plot of Fractional GC vs. Average Read Depth for the 6 largest scaffolds from the FW106 metagenome. All of these represent species of the dominant  $\gamma$ -proteobacterial species (FRC  $\gamma$  group I). The low read depth of scaffold 222 (orange) suggests that it may represent a separate  $\gamma$ -proteobacterial phylogeny.

## Phylogeny of FW106 Community

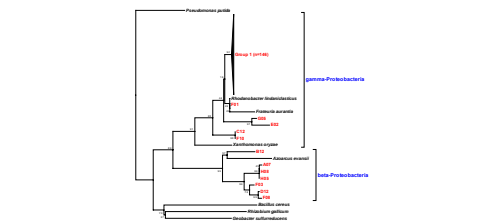


Fig. 3. Analysis of 16S rDNA from the metagenome sequence revealed 10 operational taxonomic units (OTUs) at a 99.6% similarity cutoff (2 mismatches). Phylogeny of the 16S sequences suggests that the community is dominated by a single  $\gamma$ -proteobacterial species similar to *Frateruia*.  $\beta$ -proteobacterial species were observed at lower proportions.

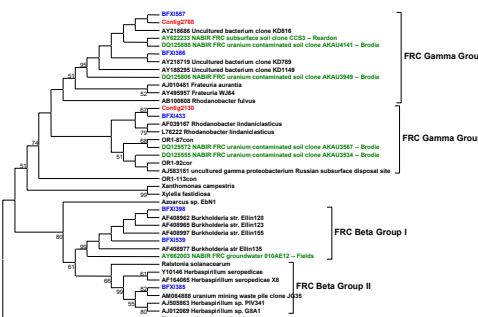


Fig. 4. Phylogeny of the predicted 16S rRNA sequences obtained from the assembled metagenome sequence. The phylogeny defines the four primary phylogenies of the community and was used for preliminary binning of the remaining contigs. Taxon labels are colored as follows: Red, genes predicted from the metagenome (2 gene fragments mapping to FRC Gamma Group I and FRC Beta Group I are not shown); Blue and Green, genes experimentally isolated from locations within the FRC.

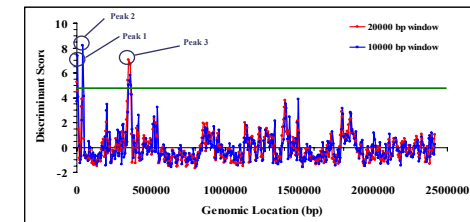
## Sampling of Genes from the Metagenome

- Energy Metabolism**
  - 30+ Cytochromes (23 c-type)
  - 3 NiFe Hydrogenase complexes
  - Assimilatory and Dissimilatory Nitrate Reductase
  - Formate Dehydrogenase
  - Multiple Cytochrome Oxidase Complexes
  - Hexose and Pentose Transporters, including multiple PTS Transport Complexes
- Stress Response**
  - RecA
  - CRP/FNR
  - Fur
  - Sigma Factors (including  $\sigma^{32}$  and  $\sigma^{38}$ )
  - 30+ Sensor Histidine Kinases
  - 40+ cation efflux/multidrug resistance pumps
  - 7 Na<sup>+</sup>/H<sup>+</sup> Antiporters
  - 9 Heat Shock Response Genes

- Metal Resistance Genes**
  - Arsenate Reductase Operon
  - Mercuric Reductase Operon
  - Copper Resistance Genes
  - CoZnCd Transport Genes
  - Chromate Transport Genes
  - Misc. Heavy Metal Resistance Genes

Genes identified from the metagenome are consistent with those predicted based on the geochemistry of the site, i.e. metal resistance genes, denitrification genes, stress response systems.

## Evidence for Lateral Gene Transfer



- Peak 1**
  - ABC-type multidrug transport operon
  - Histidine kinase/CheY-like response regulator pair
  - Transposase
  - Hypotheticals
- Peak 2**
  - RecF
  - CheY-like receiver
  - Hypotheticals
- Peak 3**
  - Multiple alcohol dehydrogenase genes
  - Amino acid transport/metabolism genes
  - Transposase
  - Hypotheticals

Fig. 5. Identification of putative genomic islands (GI) in the major scaffold 219 (FRC  $\gamma$  Group I, 2.4 Mb). GIs were determined using iterative discriminant analysis as implemented in Tu and Ding, 2003 based on GC difference, genome signature contrasts and codon usage bias. Analyses were conducted using sliding windows of 10000 (blue) and 20000 (red) bp. Peaks exceeding a discriminant score > 3.5 were considered to be potential GIs. CDS sequences in the regions of the predicted GIs are listed below the figure.

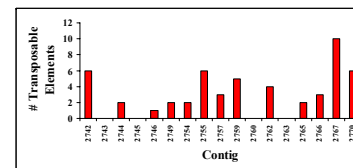


Fig. 6. Distribution of transposase, integrase, insertion elements and phage recombinase genes within the contig of the primary scaffold 219. 52 total transposable elements were identified for this scaffold.

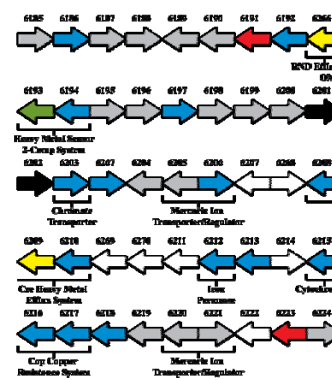


Fig. 6. The genomic environment of contig 2766 of scaffold 219 suggests evidence of lateral gene transfer of metal resistance genes into the  $\gamma$  groups from the  $\beta$  groups. ORFs are colored as follows: grey,  $\gamma$ -proteobacteria; blue,  $\beta$ -proteobacteria; yellow,  $\delta$ -proteobacteria; green,  $\alpha$ -proteobacteria; black, non-proteobacteria; white, unknown; red, transposase/integrase. ORFs are not to scale. Assignments were based on best Blast hit from the automatic annotation.

## Evolution of the FW106 Community

|                                     |               | Fixed | Polymorphic |
|-------------------------------------|---------------|-------|-------------|
| RecA                                | Nonsynonymous | 11    | 1           |
|                                     | Synonymous    | 49    | 0           |
| GroEL                               | Nonsynonymous | 4     | 1           |
|                                     | Synonymous    | 29    | 0           |
| Heavy Metal Sensor Histidine Kinase | Nonsynonymous | 14    | 1           |
|                                     | Synonymous    | 30    | 1           |
| Heavy Metal Efflux OMP, CzcC        | Nonsynonymous | 56    | 6           |
|                                     | Synonymous    | 94    | 3           |

Fig. 5. A sampling of the degree of polymorphism of genes located on the primary scaffold 219. ORF nucleotide sequences were compared to the metagenome read library using BLAST. Returned reads were clustered based on a 90% sequence identity cutoff and the resulting clusters were used to identify polymorphic sites. Sites were only considered to be polymorphic if at least two reads showed the same nucleotide change.

Based on visual analysis of the contig sequences using Consed and detailed analysis of selected individual genes, it appears that the level of detectable polymorphism in the sample is very low. The abnormal level of polymorphism in the examined *czcC* gene suggests a possible target of positive selection. The accumulated evidence suggests that the community has undergone multiple selective sweeps over the past 50 years that have significantly reduced the strain diversity of the community. The major adaptive events responsible for this phenomenon are likely the result of beneficial lateral transfer of genes between community members followed by/in conjunction with positive selection at the genetic level.

## Summary

- Groundwater community is dominated by *Frateruia*-like  $\gamma$ -proteobacteria
- Preliminary analysis suggests low strain diversity within the established phylogenies, possibly as the result of selective sweeps within the population
- Denitrifying  $\beta$ -proteobacteria may serve as keystone species
- Metabolic and regulatory genes predicted to exist based on geochemistry were identified, suggesting that the community has undergone adaptive evolution in response to contamination of the environment
- Lateral gene transfer between community members followed by positive genetic selection is proposed as the major mechanism by which the community has adapted to the environmental pressures imposed on it

## Future Work

- Metagenome-wide analysis of polymorphism
- Metagenome-wide survey of positive selection
- Metagenomic analysis of microbial communities from additional wells, including from the uncontaminated background site (proposed)
- Genomic sequencing of isolates of the major groups to test hypotheses on lateral gene transfer, positive selection, etc.

## ACKNOWLEDGEMENT

ESPP is part of the Virtual Institute for Microbial Stress and Survival supported by the U. S. Department of Energy, Office of Science, Office of Biological and Environmental Research, Genomics Program:GTL through contract DE-AC02-05CH11231 between Lawrence Berkeley National Laboratory and the U. S. Department of Energy.